# Automata & Complexity Theory
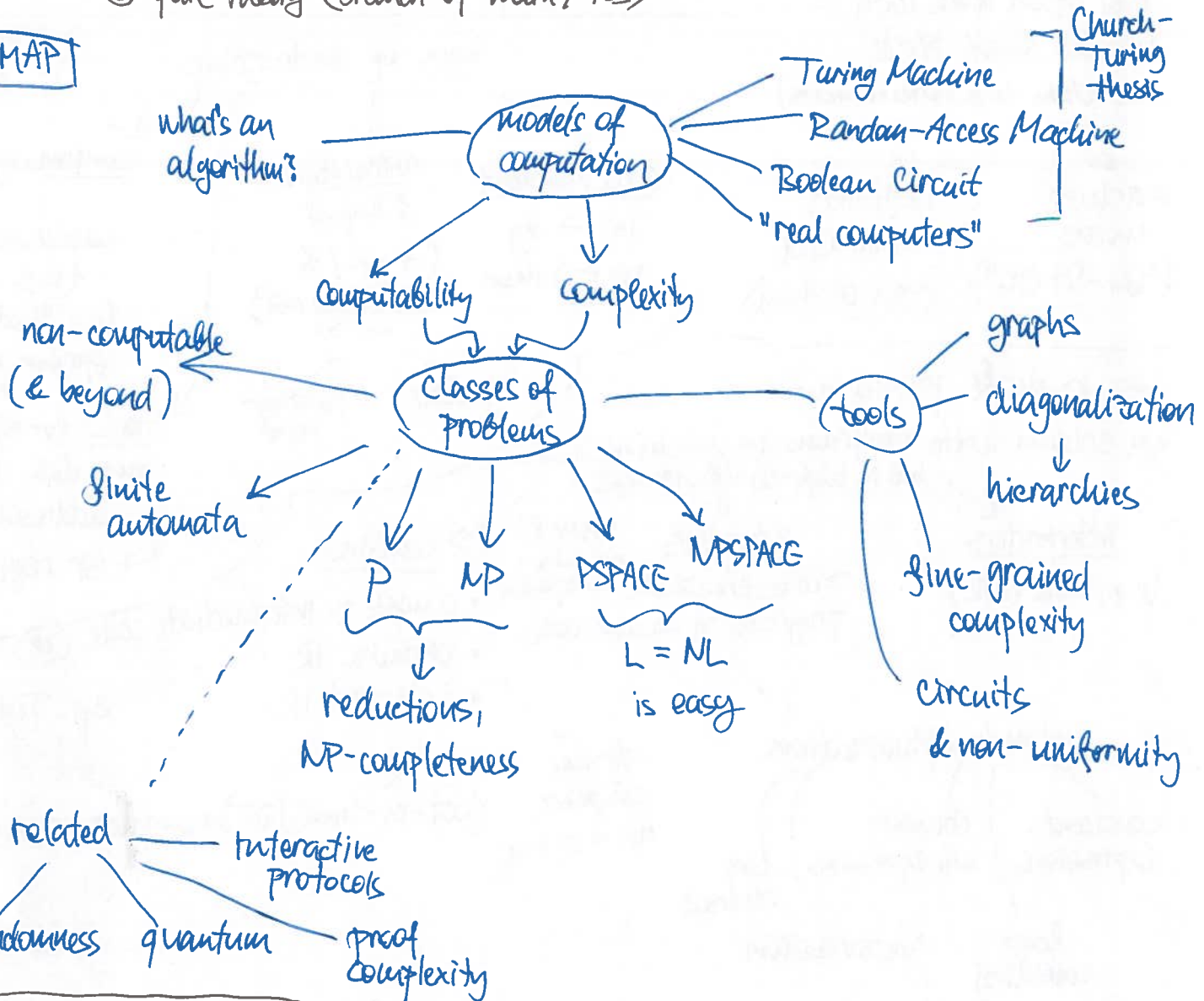
goal: build theory of computation & hardness of problems

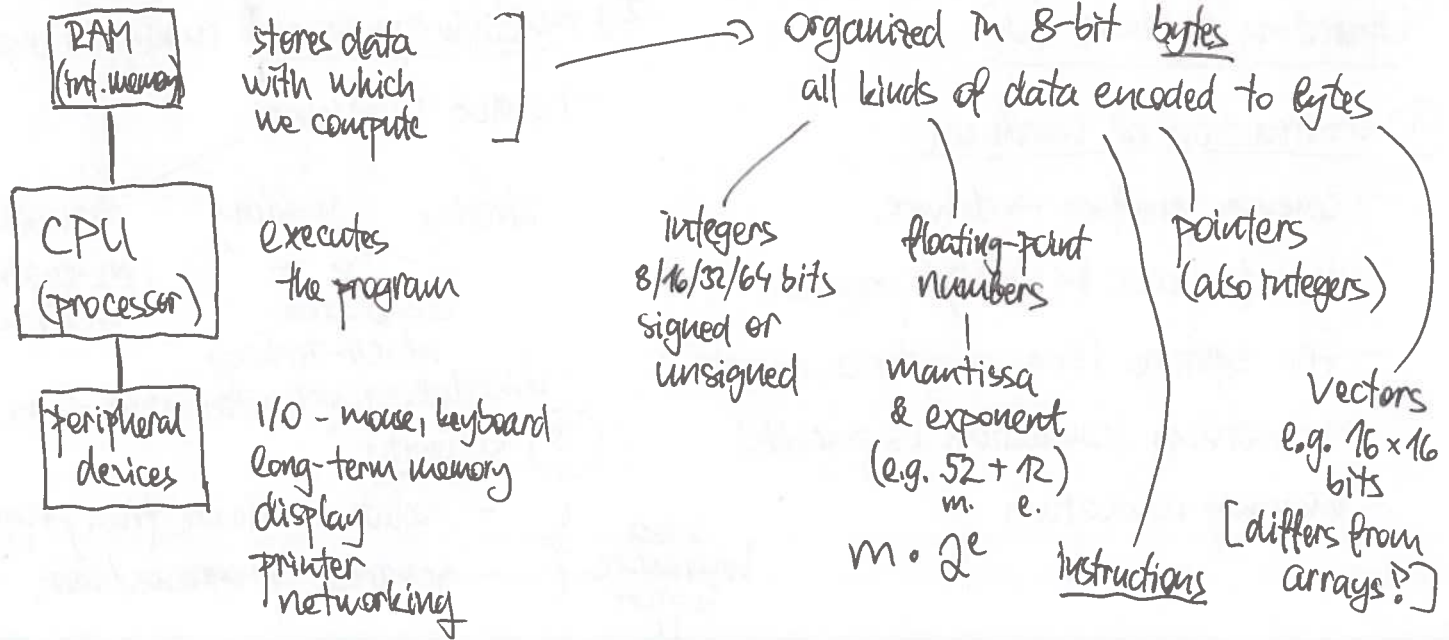two views: ① an applied theory to help us use machines more efficiently
② pure theory (branch of math / TCS)

**ROAD MAP**

what's an algorithm? → **models of computation**
- Turing Machine
- Random-Access Machine
- Boolean Circuit
- "real computers"

] Church-Turing thesis

models of computation → Computability, Complexity → **Classes of problems**

Computability → non-computable (& beyond)

**Classes of problems** →
- finite automata
- P, NP, PSPACE, NPSPACE
- related

P, NP → reductions, NP-completeness

PSPACE → $L = NL$ is easy

related → interactive protocols, randomness, quantum, proof complexity

**Classes of problems** → **tools** →
- graphs
- diagonalization ↓ hierarchies
- fine-grained complexity
- circuits & non-uniformity

---

**PHYSICAL COMPUTERS** & their architecture (typical case in 2022, just sketching)

**RAM** (int. memory) — stores data with which we compute ] → organized in 8-bit **bytes**
all kinds of data encoded to bytes

**CPU** (processor) — executes the program

**peripheral devices** — I/O - mouse, keyboard
long-term memory
display
printer
networking

integers
8/16/32/64 bits
signed or unsigned

floating-point numbers
↓
mantissa & exponent
(e.g. 52 + 12)
   m.   e.

$m \cdot 2^e$

pointers (also integers)

vectors
e.g. 16 × 16 bits
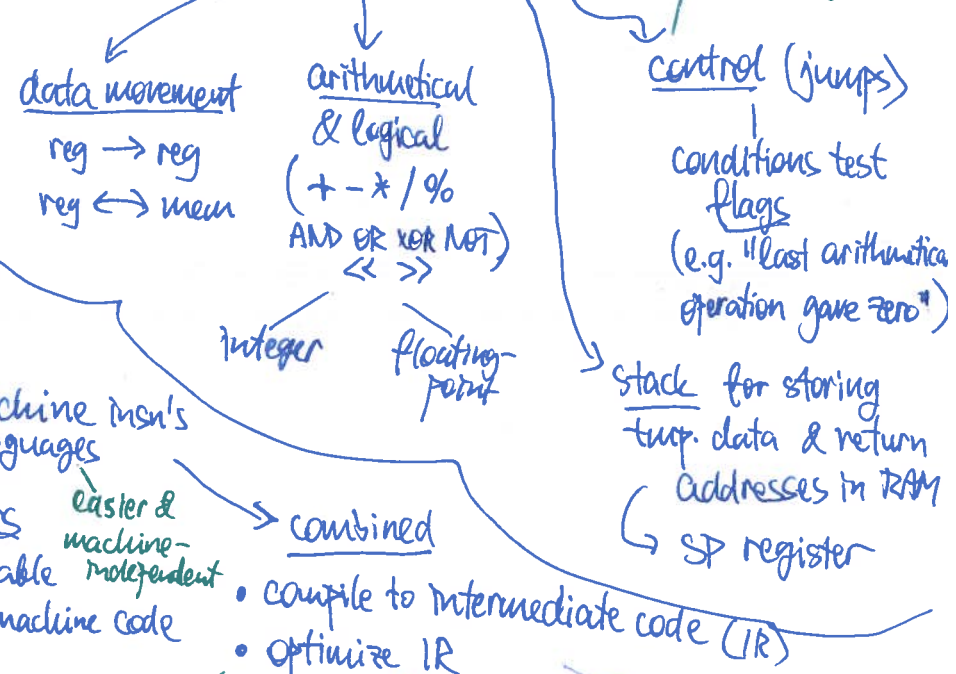[differs from arrays!?]

instructions

machine instructions
(program)

- stored in the same memory as data (Von Neumann architecture)
- can modify itself
- just another interpretation of bytes

↓

they often work with several simple pieces of data (e.g., 64b numbers)

↓ ↓

machine words ("64-bit CPU")

registers inside CPU (~10s of them)

---

types of instructions

→ program flow, PC register

**data movement**
reg → reg
reg ↔ mem

**arithmetical & logical**
( + - * / %
AND OR XOR NOT)
<< >>

integer    floating-point

**control** (jumps)
↓
conditions test flags
(e.g. "last arithmetical operation gave zero")

stack for storing temp. data & return addresses in RAM
↳ SP register

---

see example program ...

We seldom write programs in machine insn's but in high-level languages

**interpreters**
(e.g., UNIX shell)

**compilers**
produce executable programs in machine code

easier & machine-independent

→ **combined**
• compile to intermediate code (IR)
• optimize IR
• interpret IR

e.g. Python

↓

automatic optimization

constant expressions | common sub-expressions | loop reversal

loop unrolling    vectorization

& many others

typical compiler:
HLL → IR → MC

just-in-time (JIT) compilers — e.g. Java

---

**Operating systems (OS)**

① **abstraction of hardware**

- common interface → drivers
- manage access by multiple programs
- file systems (files, directories, mounts...)
- networking (connections vs. packets)
- memory allocation

② **multiple processes** "running at once"

- context switching

sleeping    yielding    timeslices
(pre-emptive multitasking)

cooperative multi-tasking

- scheduling, priorities, real-time

③ **security**

need hardware support ⎰ - isolate hardware from programs
                      ⎱ - separate programs/users

HW features for security — privilege levels (supervisor/user mode) → system call instructions ③

Virtual memory management

virtual memory (4kB pages) → MMU → physical memory (same-size pages) → page fault exception

mapping of virt. to phys. page numbers —— includes access rights ← read R / write W / execute code X / supervisor only (or switching tables on mode changes...)

Uses: — protection of processes (private memory) ... RW(X) for 1 process
— shared memory ... RW(X) in multiple processes
— shared library ... R in multiple processes
— lazy allocation ... read-only shared zeroes, copy on write
— fork ... using copy-on-write mapping
— swapping ... store seldom-used pages to disk, read back on access

• caching — RAM is slow (CPU executes ~ $10^9$ instructions/second, RAM latency is tens of ns)
— idea: small, very fast memory inside the CPU which remembers frequently used data ] called a cache | can be better only for small memory (speed of light &c.)

— caches 64B chunks of data (cache lines)

— strategy: — write-through vs. write-back
— when cache fills up: evict least-recently used item (LRU)

— real caches have limited associativity → cache aliasing
— multiple levels of caches
— example: accessing a matrix row by row vs. column by column
↳ sequential in memory   ↳ every access is a cache miss → very slow

→ modelling of caches, cache-oblivious algorithms (not at this lecture)
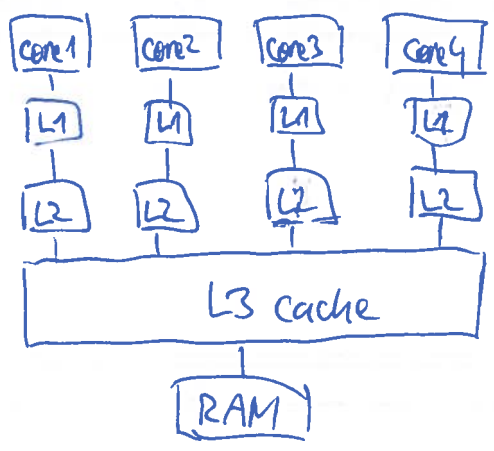
• improving execution of instructions
— CPU works in cycles, historically 1 instruction took multiple cycles   e.g.:
— pipelining  [fetch|load|comp|store] ins.1
                [fetch|load|comp|store] ins.2
                T1  T2  T3  T4  T5
all units of CPU always busy
problems: dependencies & (conditional) jumps
① fetch & decode
② load operands
③ compute result
④ store result

— superscalar CPU: multiple units for different types of instructions, can run in parallel
→ scheduling instructions to units
all this is transparent to SW
(well, almost: Meltdown & other bugs)

— jump prediction

- **multiple processors** sharing memory (SMP = Symmetric Multi-Processing)
  - OS schedules processes on processors — real parallelism
  - hard to get right: locking in SW, cache coherency protocols in HW
- **multi-core processors**: SMP on a single chip

For example:

| Core1 | Core2 | Core3 | Core4 |

| L1 | L1 | L1 | L1 |

| L2 | L2 | L2 | L2 |

L3 cache

RAM

typical sizes

32 KB code + 32 KB data

256 KB unified

8 MB unified

16 GB

- **multi-threaded cores**: two cores sharing their execution units & caches
  - unclear benefits (can even make things worse!)
- **virtual machines**: simulating a whole machine within a process
  - including supervisor mode → the VM can run its own OS
  - including virtual peripherals
  - CPUs have special support for VMs (e.g., nested paging in MMU)

## Relationship with theory

- will ignore most machine-dependent constants
- concentrate on asymptotics ⇒ all machines (roughly) equal
- use simple mathematical machines instead
- I/O and caches need special treatment

} rest of the semester